

1c912 U.S. PTO  
09/13/00

09-15-00

A

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

Inventorship..... Rui et al.  
Applicant..... Microsoft Corporation  
Attorney's Docket No. .... MS1-416US  
Title: Annotating Programs for Automatic Summary Generation

1c662 U.S. PTO  
09/660529  
09/13/00

TRANSMITTAL LETTER AND CERTIFICATE OF MAILING

To: Commissioner of Patents and Trademarks,  
Washington, D.C. 20231

From: Allan T. Sponseller (Tel. 509-324-9256; Fax 509-323-8979)  
Lee & Hayes, PLLC  
421 W. Riverside Avenue, Suite 500  
Spokane, WA 99201

The following enumerated items accompany this transmittal letter and are being submitted for the matter identified in the above caption.

1. Specification—title page, plus 48 pages, including 65 claims and Abstract
2. Transmittal letter including Certificate of Express Mailing
3. 7 Sheets Formal Drawings (Figs. 1-7)
4. Return Post Card

Large Entity Status [x]                      Small Entity Status [ ]

Date: Sept. 13, 2000

By: [Signature]  
Allan T. Sponseller  
Reg. No. 38,318

CERTIFICATE OF MAILING

I hereby certify that the items listed above as enclosed are being deposited with the U.S. Postal Service as either first class mail, or Express Mail if the blank for Express Mail No. is completed below, in an envelope addressed to The Commissioner of Patents and Trademarks, Washington, D.C. 20231, on the below-indicated date. Any Express Mail No. has also been marked on the listed items.

Express Mail No. (if applicable) EL624352652

Date: 9-13-00

By: [Signature]  
Lori A. Vierra

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

**Annotating Programs for Automatic Summary  
Generation**

Inventor(s):

Yong Rui

Anoop Gupta

Alejandro Acero

1 **RELATED APPLICATIONS**

2 This application claims the benefit of U.S. Provisional Application No.  
3 60/153,730, filed September 13, 1999, entitled "MPEG-7 Enhanced Multimedia  
4 Access" to Yong Rui, Jonathan Grudin, Anoop Gupta, and Liwei He, which is  
5 hereby incorporated by reference.

6  
7 **TECHNICAL FIELD**

8 This invention relates to audio/video programming and rendering thereof,  
9 and more particularly to annotating programs for automatic summary generation.

10  
11 **BACKGROUND OF THE INVENTION**

12 Watching television has become a common activity for many people,  
13 allowing people to receive important information (e.g., news broadcasts, weather  
14 forecasts, etc.) as well as simply be entertained. While the quality of televisions  
15 on which programs are rendered has improved, so too have a wide variety of  
16 devices been developed and made commercially available that further enhance the  
17 television viewing experience. Examples of such devices include Internet  
18 appliances that allow viewers to "surf" the Internet while watching a television  
19 program, recording devices (either analog or digital) that allow a program to be  
20 recorded and viewed at a later time, etc.

21 Despite these advances and various devices, mechanisms for watching  
22 television programs are still limited to two general categories: (1) watching the  
23 program "live" as it is broadcast, or (2) recording the program for later viewing.  
24 Each of these mechanisms, however, limits viewers to watching their programs in  
25 the same manner as they were was broadcast (although possibly time-delayed).

Often times, however, people do not have sufficient time to watch the entirety of a recorded television program. By way of example, a sporting event such as a baseball game may take 2 or 2½ hours, but a viewer may only have ½ hour that he or she can spend watching the recorded game. Currently, the only way for the viewer to watch such a game is for the viewer to randomly select portions of the game to watch (e.g., using fast forward and/or rewind buttons), or alternatively use a "fast forward" option to play the video portion of the recorded game back at a higher speed than that at which it was recorded (although no audio can be heard). Such solutions, however, have significant drawbacks because it is extremely difficult for the viewer to know or identify which portions of the game are the most important for him or her to watch. For example, the baseball game may have only a handful of portions that are exciting, with the rest being uninteresting and not exciting.

The invention described below addresses these disadvantages, providing for annotating of programs for automatic summary generation.

## **SUMMARY OF THE INVENTION**

Annotating programs for automatic summary generation is described herein.

In accordance with one aspect, audio/video programming content is made available to a receiver from a content provider, and meta data is made available to the receiver from a meta data provider. The content provider and meta data provider may be the same or different devices. The meta data corresponds to the programming content, and identifies, for each of multiple portions of the programming content, an indicator of a likelihood that the portion is an exciting

1 portion of the content. The meta data can be used, for example, to allow  
2 summaries of the programming content to be generated by selecting the portions  
3 having the highest likelihoods of being exciting portions.

4 According to another aspect, exciting portions of a sporting event are  
5 automatically identified based on sports-specific events and sports-generic events.  
6 The audio data of the sporting event is analyzed to identify sports-specific events  
7 (such as baseball hits if the sporting event is a baseball program) as well as sports-  
8 generic events (such as excited speech from an announcer). These sports-specific  
9 and sports-generic events are used together to identify the exciting portions of the  
10 sporting event.

11 According to another aspect, exciting segments of a baseball program are  
12 automatically identified. Various features are extracted from the audio data of the  
13 baseball program and selected features are input to an excited speech classification  
14 subsystem and a baseball hit detection subsystem. The excited speech  
15 classification subsystem identifies probabilities that segments of the audio data  
16 contain excited speech (e.g., from an announcer). The baseball hit detection  
17 subsystem identifies probabilities that multiple-frame groupings of the audio data  
18 include baseball hits. These two sets of probabilities are input to a probabilistic  
19 fusion subsystem that determines, based on both probabilities, a likelihood that  
20 each of the segments is an exciting portion of the baseball program. These  
21 probabilities can then be used, for example, to generate a summary of the baseball  
22 program.

## **BRIEF DESCRIPTION OF THE DRAWINGS**

The present invention is illustrated by way of example and not limitation in the figures of the accompanying drawings. The same numbers are used throughout the figures to reference like components and/or features.

Fig. 1 shows a programming distribution and viewing system in accordance with one embodiment of the invention;

Fig. 2 illustrates an example of a suitable operating environment in which the invention may be implemented;

Fig. 3 illustrates an exemplary programming content delivery architecture in accordance with certain embodiments of the invention;

Fig. 4 illustrates an exemplary automatic summary generation process in accordance with certain embodiments of the invention;

Fig. 5 illustrates part of an exemplary audio clip and portions from which features are extracted;

Fig. 6 illustrates exemplary baseball hit templates that may be used in accordance with certain embodiments of the invention; and

Fig. 7 is a flowchart illustrating an exemplary process for rendering a program summary to a user in accordance with certain embodiments of the invention.

## **DETAILED DESCRIPTION**

### **General System**

Fig. 1 shows a programming distribution and viewing system 100 in accordance with one embodiment of the invention. System 100 includes a video and audio rendering system 102 having a display device including a viewing area

1 104. Video and audio rendering system 102 represents any of a wide variety of  
2 devices for playing video and audio content, such as a traditional television  
3 receiver, a personal computer, etc. Receiver 106 is connected to receive and  
4 render content from multiple different programming sources. Although illustrated  
5 as separate components, rendering system 102 may be combined with receiver 106  
6 into a single component (e.g., a personal computer or television). Receiver 106  
7 may also be capable of storing content locally, in either analog or digital format  
8 (e.g., on magnetic tapes, a hard disk drive, optical disks, etc.).

9 While audio and video have traditionally been transmitted using analog  
10 formats over the airwaves, current and proposed technology allows multimedia  
11 content transmission over a wider range of network types, including digital  
12 formats over the airwaves, different types of cable and satellite systems  
13 (employing both analog and digital transmission formats), wired or wireless  
14 networks such as the Internet, etc.

15 Fig. 1 shows several different physical sources of programming, including a  
16 terrestrial television broadcasting system 108 which can broadcast analog or  
17 digital signals that are received by antenna 110; a satellite broadcasting system 112  
18 which can transmit analog or digital signals that are received by satellite dish 114;  
19 a cable signal transmitter 116 which can transmit analog or digital signals that are  
20 received via cable 118; and an Internet provider 120 which can transmit digital  
21 signals that are received by modem 122 via the Internet (and/or other network)  
22 124. Both analog and digital signals can include programming made up of audio,  
23 video, and/or other data. Additionally, a program may have different components  
24 received from different programming sources, such as audio and video data from  
25 cable transmitter 116 but data from Internet provider 120. Other programming

1 sources might be used in different situations, including interactive television  
2 systems.

3 As described in more detail below, programming content made available to  
4 system 102 includes audio and video programs as well as meta data corresponding  
5 to the programs. The meta data is used to identify portions of the program that are  
6 believed to be exciting portions, as well as how exciting these portions are  
7 believed to be relative to one another. The meta data can be used to generate  
8 summaries for the programs, allowing the user to view only the portions of the  
9 program that are determined to be the most exciting.

### 10 11 **Exemplary Operating Environment**

12 Fig. 2 illustrates an example of a suitable operating environment in which  
13 the invention may be implemented. The illustrated operating environment is only  
14 one example of a suitable operating environment and is not intended to suggest  
15 any limitation as to the scope of use or functionality of the invention. Other well  
16 known computing systems, environments, and/or configurations that may be  
17 suitable for use with the invention include, but are not limited to, personal  
18 computers, server computers, hand-held or laptop devices, multiprocessor systems,  
19 microprocessor-based systems, programmable consumer electronics (e.g., digital  
20 video recorders), gaming consoles, cellular telephones, network PCs,  
21 minicomputers, mainframe computers, distributed computing environments that  
22 include any of the above systems or devices, and the like.

23 Alternatively, the invention may be implemented in hardware or a  
24 combination of hardware, software, and/or firmware. For example, one or more  
25



1 application specific integrated circuits (ASICs) could be designed or programmed  
2 to carry out the invention.

3 Fig. 2 shows a general example of a computer 142 that can be used in  
4 accordance with the invention. Computer 142 is shown as an example of a  
5 computer that can perform the functions of receiver 106 of Fig. 1, or of one of the  
6 programming sources of Fig. 1 (e.g., Internet provider 120). Computer 142  
7 includes one or more processors or processing units 144, a system memory 146,  
8 and a bus 148 that couples various system components including the system  
9 memory 146 to processors 144.

10 The bus 148 represents one or more of any of several types of bus  
11 structures, including a memory bus or memory controller, a peripheral bus, an  
12 accelerated graphics port, and a processor or local bus using any of a variety of  
13 bus architectures. The system memory 146 includes read only memory (ROM)  
14 150 and random access memory (RAM) 152. A basic input/output system (BIOS)  
15 154, containing the basic routines that help to transfer information between  
16 elements within computer 142, such as during start-up, is stored in ROM 150.  
17 Computer 142 further includes a hard disk drive 156 for reading from and writing  
18 to a hard disk, not shown, connected to bus 148 via a hard disk drive interface 157  
19 (e.g., a SCSI, ATA, or other type of interface); a magnetic disk drive 158 for  
20 reading from and writing to a removable magnetic disk 160, connected to bus 148  
21 via a magnetic disk drive interface 161; and an optical disk drive 162 for reading  
22 from and/or writing to a removable optical disk 164 such as a CD ROM, DVD, or  
23 other optical media, connected to bus 148 via an optical drive interface 165. The  
24 drives and their associated computer-readable media provide nonvolatile storage  
25 of computer readable instructions, data structures, program modules and other data

1 for computer 142. Although the exemplary environment described herein employs  
2 a hard disk, a removable magnetic disk 160 and a removable optical disk 164, it  
3 will be appreciated by those skilled in the art that other types of computer readable  
4 media which can store data that is accessible by a computer, such as magnetic  
5 cassettes, flash memory cards, random access memories (RAMs), read only  
6 memories (ROM), and the like, may also be used in the exemplary operating  
7 environment.

8 A number of program modules may be stored on the hard disk, magnetic  
9 disk 160, optical disk 164, ROM 150, or RAM 152, including an operating system  
10 170, one or more application programs 172, other program modules 174, and  
11 program data 176. A user may enter commands and information into computer  
12 142 through input devices such as keyboard 178 and pointing device 180. Other  
13 input devices (not shown) may include a microphone, joystick, game pad, satellite  
14 dish, scanner, or the like. These and other input devices are connected to the  
15 processing unit 144 through an interface 168 that is coupled to the system bus  
16 (e.g., a serial port interface, a parallel port interface, a universal serial bus (USB)  
17 interface, etc.). A monitor 184 or other type of display device is also connected to  
18 the system bus 148 via an interface, such as a video adapter 186. In addition to the  
19 monitor, personal computers typically include other peripheral output devices (not  
20 shown) such as speakers and printers.

21 Computer 142 operates in a networked environment using logical  
22 connections to one or more remote computers, such as a remote computer 188.  
23 The remote computer 188 may be another personal computer, a server, a router, a  
24 network PC, a peer device or other common network node, and typically includes  
25 many or all of the elements described above relative to computer 142, although

only a memory storage device 190 has been illustrated in Fig. 2. The logical connections depicted in Fig. 2 include a local area network (LAN) 192 and a wide area network (WAN) 194. Such networking environments are commonplace in offices, enterprise-wide computer networks, intranets, and the Internet. In certain embodiments of the invention, computer 142 executes an Internet Web browser program (which may optionally be integrated into the operating system 170) such as the “Internet Explorer” Web browser manufactured and distributed by Microsoft Corporation of Redmond, Washington.

When used in a LAN networking environment, computer 142 is connected to the local network 192 through a network interface or adapter 196. When used in a WAN networking environment, computer 142 typically includes a modem 198 or other means for establishing communications over the wide area network 194, such as the Internet. The modem 198, which may be internal or external, is connected to the system bus 148 via a serial port interface 168. In a networked environment, program modules depicted relative to the personal computer 142, or portions thereof, may be stored in the remote memory storage device. It will be appreciated that the network connections shown are exemplary and other means of establishing a communications link between the computers may be used.

Computer 142 also includes a broadcast tuner 200. Broadcast tuner 200 receives broadcast signals either directly (e.g., analog or digital cable transmissions fed directly into tuner 200) or via a reception device (e.g., via antenna 110 or satellite dish 114 of Fig. 1).

Computer 142 typically includes at least some form of computer readable media. Computer readable media can be any available media that can be accessed by computer 142. By way of example, and not limitation, computer readable

1 media may comprise computer storage media and communication media.  
2 Computer storage media includes volatile and nonvolatile, removable and non-  
3 removable media implemented in any method or technology for storage of  
4 information such as computer readable instructions, data structures, program  
5 modules or other data. Computer storage media includes, but is not limited to,  
6 RAM, ROM, EEPROM, flash memory or other memory technology, CD-ROM,  
7 digital versatile disks (DVD) or other optical storage, magnetic cassettes, magnetic  
8 tape, magnetic disk storage or other magnetic storage devices, or any other media  
9 which can be used to store the desired information and which can be accessed by  
10 computer 142. Communication media typically embodies computer readable  
11 instructions, data structures, program modules or other data in a modulated data  
12 signal such as a carrier wave or other transport mechanism and includes any  
13 information delivery media. The term "modulated data signal" means a signal that  
14 has one or more of its characteristics set or changed in such a manner as to encode  
15 information in the signal. By way of example, and not limitation, communication  
16 media includes wired media such as wired network or direct-wired connection,  
17 and wireless media such as acoustic, RF, infrared and other wireless media.  
18 Combinations of any of the above should also be included within the scope of  
19 computer readable media.

20 The invention has been described in part in the general context of  
21 computer-executable instructions, such as program modules, executed by one or  
22 more computers or other devices. Generally, program modules include routines,  
23 programs, objects, components, data structures, etc. that perform particular tasks  
24 or implement particular abstract data types. Typically the functionality of the  
25

1 program modules may be combined or distributed as desired in various  
2 embodiments.

3 For purposes of illustration, programs and other executable program  
4 components such as the operating system are illustrated herein as discrete blocks,  
5 although it is recognized that such programs and components reside at various  
6 times in different storage components of the computer, and are executed by the  
7 data processor(s) of the computer.

### 8 9 **Content Delivery Architecture**

10 Fig. 3 illustrates an exemplary programming content delivery architecture  
11 in accordance with certain embodiments of the invention. A client 220 receives  
12 programming content including both audio/video data 222 and meta data 224 that  
13 corresponds to the audio/video data 222. In the illustrated example, an  
14 audio/video data provider 226 is the source of audio/video data 222 and a meta  
15 data provider 228 is the source of meta data 224. Alternatively, meta data 224 and  
16 audio/video data 222 may be provided by the same source, or alternatively three or  
17 more different sources.

18 The data 222 and 224 can be made available by providers 226 and 228 in  
19 any of a wide variety of formats. In one implementation, data 222 and 224 are  
20 formatted in accordance with the MPEG-7 (Moving Pictures Expert Group)  
21 format. The MPEG-7 format standardizes a set of Descriptors (Ds) that can be  
22 used to describe various types of multimedia content, as well as a set of  
23 Description Schemes (DSs) to specify the structure of the Ds and their  
24 relationship. In MPEG-7, the audio and video data 222 are each described as one  
25 or more Descriptors, and the meta data 224 is described as a Description Scheme.

Client 220 is illustrated as separate from providers 226 and 228. This separation can be small (e.g., across a LAN) or large (e.g., a remote server located in another city or state). Alternatively, data 222 and/or 224 may be stored locally by client 220, either on another device such as an analog or digital video recorder (not shown) coupled to client 220 or within client 220 (e.g., on a hard disk drive).

A wide variety of meta data 224 can be associated with a program. In the discussions below, meta data 224 is described as being "excited segment probabilities" which identify particular segments of the program and a corresponding probability or likelihood that each segment is an "exciting" segment. An exciting segment is a segment of the program believed to be typically considered exciting to viewers. By way of example, baseball hits are believed to be typically considered exciting segments of a baseball program.

The excited segment probabilities in meta data 224 can be generated in any of a variety of manners. In one implementation, the excited segment probabilities

are generated manually (e.g., by a producer or other individual(s) watching the program and identifying the exciting segments and assigning the corresponding probabilities). In another implementation, the excited segment probabilities are generated automatically by a process described in more detail below. Additionally, the excited segment probabilities can be generated after the fact (e.g., after a baseball game is over and its entirety is available on a recording medium), or alternatively on the fly (e.g., a baseball game may be monitored and probabilities generated as the game is played).

### **Automatic Summary Generation**

The automatic summary generation process described below refers to sports-generic and sports-specific events, and refers specifically to the example of a baseball program. Alternatively, summaries can be automatically generated in an analogous manner for other programs, including other sporting events.

The automatic summary generation process analyzes the audio data of the baseball program and attempts to identify segments that include speech, and of those segments which can be identified as being "excited" speech (e.g., the excitement in an announcer's voice). Additionally, based on the audio data segments that include baseball hits are also identified. These excited speech segments and baseball hit segments are then used to determine, for each of the excited speech segments, a probability that the segment is truly an exciting segment of the program. Given these probabilities, a summary of the program can be generated.

Fig. 4 illustrates an exemplary automatic summary generation process in accordance with certain embodiments of the invention. The generation process

1 begins with the raw audio data 250 (also referred to as a raw audio clip), such as  
2 the audio portion of data 222 of Fig. 3. The raw audio data 250 is the audio  
3 portion of the program for which the summary is being automatically generated.  
4 The audio data 250 is input to feature extractor 252 which extracts various features  
5 from portions of audio data 250. In one implementation, feature extractor 252  
6 extracts one or more of energy features, phoneme-level features, information  
7 complexity features, and prosodic features.

8 Fig. 5 illustrates part of an exemplary audio clip and portions from which  
9 features are extracted. Audio clip 258 is illustrated. Audio features are extracted  
10 from audio clip 258 using two different resolutions: a sports-specific event  
11 detection resolution used to assist in the identification of potentially exciting  
12 sports-specific events, and a sports-generic event detection resolution used to  
13 assist in the identification of potentially exciting sports-generic events. In the  
14 illustrated example, the sports-specific event detection resolution is 10  
15 milliseconds (ms), while the sports-generic event detection resolution is 0.5  
16 seconds. Alternatively, other resolutions could be used.

17 As used herein, the sports-specific event detection is based on 10 ms  
18 "frames", while the sports-generic event detection is based on 0.5 second  
19 "windows". As illustrated in Fig. 5, the 10 ms frames are non-overlapping and the  
20 0.5 second windows are non-overlapping, although the frames overlap the  
21 windows (and vice versa). Alternatively, the frames may overlap other frames,  
22 and/or the windows may overlap other windows.

23 Returning to Fig. 4, feature extractor 252 extracts different features from  
24 audio data 250 based on both frames and windows of audio data 250. Exemplary  
25 features which can be extracted by feature extractor 252 are discussed below.





1       Extractor 252 extracts phoneme-level features for each of the 10ms frames  
2 of audio data 250. For each frame, two well-known feature vectors are extracted:  
3 a Mel-frequency Cepstral coefficient (MFCC) and the first derivative of the  
4 MFCC (referred to as the delta MFCC). The MFCC is the *cosine* transform of the  
5 pitch of the frame on the "Mel-scale", which is a gradually warped linear spectrum  
6 (with coarser resolution at high frequencies).

7       Extractor 252 extracts information complexity features for each of the 10  
8 ms frames of audio data 250. For each frame, a feature vector representing the  
9 entropy (*Etr*) of the frame is extracted. For an *N*-point Fast Fourier Transform  
10 (FFT) of an audio signal *s(t)*, with *S(n)* representing the *n*th frequency's  
11 component, entropy is defined as:

$$Etr = \sum_{n=1}^N P_n \log P_n$$

12  
13  
14  
15 where:

$$P_n = \frac{|S(n)|^2}{\sum_{n=1}^N |S(n)|^2}$$

16  
17  
18  
19       Extracting feature vectors representing entropy is well-known to those  
20 skilled in the art and thus will not be discussed further except as it relates to the  
21 present invention.

22       Extractor 252 extracts prosodic features for each of the 0.5 second windows  
23 of audio data 250. For each window, a feature vector representing the pitch (*Pch*)  
24 of the window is extracted. A variety of different well-known approaches can be  
25

used in determining pitch, such as the auto-regressive model, the average magnitude difference function, the maximum *a posteriori* (MAP) approach, etc.

The pitch is also determined for each 10ms frame of the 0.5 second window. These individual frame pitches are then used to extract pitch statistics regarding the pitch of the window. Exemplary pitch statistics extracted for each 0.5 second window are illustrated in Table II.

Table II

Statistic	Description
non-zero pitch count	The number of frames in the window that have a non-zero pitch value.
maximum pitch	The highest pitch value of the frames in the window.
minimum pitch	The lowest pitch value of the frames in the window.
average pitch	The average pitch value of the frames in the window.
pitch dynamic range	The pitch range over the frames in the window (the difference between the maximum and minimum pitch values).

Selected ones of the extracted features are passed by feature extractor 252 to an excited speech classification subsystem 260 and a baseball hit detection subsystem 262. Excited speech classification subsystem 260 attempts to identify segments of the audio data that include excited speech (sports-generic events), while baseball hit detection subsystem 262 attempts to identify segments of the audio data that include baseball hits (sports-specific events). The segments identified by subsystems 260 and 262 may be of the same or alternatively different sizes (and may be varying sizes). Probabilities generated for the segments are then

Excited speech classification subsystem 260 uses a two-stage process to identify segments of excited speech. In a first stage, energy and phoneme-level features 266 from feature extractor 252 are input to a speech detector 268 that identifies windows of the audio data that include speech (speech windows 270). In the illustrated example, speech detector 268 uses both the  $E_{23}$  and the delta MFCC feature vectors. For each 0.5 second window, if the  $E_{23}$  and delta MFCC vectors each exceed corresponding thresholds, the window is identified as a speech window 270; otherwise, the window is classified as not including speech. In one implementation, the thresholds used by speech detector 268 are 2.0 for the delta MFCC feature, and  $0.07 * E_{cap}$  for the  $E_{23}$  feature (where  $E_{cap}$  is the highest  $E_{23}$  value of all the frames in the audio clip (or alternatively all of the frames in the audio clip that have been analyzed so far), although different thresholds could alternatively be used.

In alternative embodiments, speech detector 268 may use different features to classify segments as speech or not speech. By way of example, energy only may be used (e.g., the window is classified as speech only if  $E_{23}$  exceeds a threshold amount (such as  $0.2 * E_{cap}$ ). By way of another example, energy and entropy features may both be used (e.g., the window is classified as speech only if the product of  $E_{23}$  and  $E_{tr}$  exceeds a threshold amount (such as 50,000).

In the second stage, pitch and energy features 272, received from feature extractor 252, for each of the speech windows 270 are used by excited speech classifier 274 to determine a probability that each speech window 270 is excited speech. Classifier 274 then combines these probabilities to identify a probability

Excited speech classifier 274 uses six statistics regarding the energy  $E_{23}$  features and the pitch ( $Pch$ ) features extracted from each speech window 270: maximum energy, average energy, energy dynamic range, maximum pitch, average pitch, and pitch dynamic range. Classifier 274 concatenates these six statistics together to generate a feature vector (having nine elements or dimensions) and compares the feature vector to a set of training vectors (based on corresponding features of training sample data) in two different classes: an excited speech class and a non-excited speech class. The *posterior* probability of a feature vector  $X$  (for a window 270) being in a class  $C_i$ , where  $C_1$  is the class of excited speech and  $C_2$  is the class of non-excited speech, can be represented as:  $P(C_i | X)$ . The probability of error in classifying the feature vector  $X$  can be reduced by classifying the data to the class having the *posterior* probability that is the highest.

Speech classifier 274 determines the *posterior* probability  $P(C_i | X)$  using learning machines. A wide variety of different learning machines can be used to determine the *posterior* probability  $P(C_i | X)$ . Three such learning machines are described below, although other learning machines could alternatively be used.

The *posterior* probability  $P(C_i | X)$  can be determined using parametric machines, such as Bayes rule:

$$P(C_i | X) = \frac{P(C_i)p(X | C_i)}{p(X)}$$

where  $p(X)$  is the data density,  $P(C_i)$  is the prior probability, and  $p(X | C_i)$  is the conditional class density. The data density  $p(x)$  is a constant for all the classes and thus does not contribute to the decision rule. The prior probability  $P(C_i)$  can be estimated from labeled training data (e.g., excited speech and non-excited speech) in a conventional manner. The conditional class density  $p(X | C_i)$  can be calculated in a variety of different manners, such as the Gaussian (Normal) distribution  $N(\mu, \sigma)$ . The  $\mu$  parameter (mean) and the  $\sigma$  parameter (standard deviation) can be determined using the well-known Maximum Likelihood Estimation (MLE):

$$\mu = \frac{1}{n} \sum_{k=1}^n X_k$$

$$\sigma^2 = \frac{1}{n} \sum_{k=1}^n (X_k - \mu)^2$$

where  $n$  is the number of training samples and  $X$  represents the training samples.

Another type of machines that can be used to determine the *posterior* probability  $P(C_i | X)$  are non-parametric machines. The K nearest neighbor technique is an example of such a machine. Using the K nearest neighbor technique:

$$P(C_i | X) = \frac{\frac{K_i}{nV}}{\sum_i \frac{K_i}{nV}} = \frac{K_i}{K}$$

where  $V$  is the volume around feature vector  $X$ ,  $V$  covers  $K$  labeled (training) samples, and  $K_i$  is the number of samples in class  $C_i$ .

Another type of machines that can be used to determine the *posterior* probability  $P(C_i | X)$  are semi-parametric machines, which combine the advantages of non-parametric and parametric machines. Examples of such semi-parametric machines include Gaussian mixture models, neural networks, and support vector machines (SVMs).

Any of a wide variety of well-known training methods can be used to train the SVM. After the SVM is trained, a sigmoid function is trained to map the SVM outputs into *posterior* probabilities. The *posterior* probability  $P(C_i | X)$  can then be determined as follows:

$$P(C_i | X) = \frac{1}{1 + \exp(AX + B)}$$

where  $A$  and  $B$  are the parameters of the sigmoid function. The parameters  $A$  and  $B$  are determined by reducing the negative log likelihood of training data  $(f_i, t_i)$ , which is a cross-entropy error function:

$$\min - \sum_i t_i \log(p_i) + (1 - t_i) \log(1 - p_i)$$

where

$$p_i = \frac{1}{1 + \exp(Af_i + B)}$$

The cross-entropy error function minimization can be performed using any number of conventional optimization processes. The training data  $(f_i, t_i)$  can be the same training data used to train the SVM, or other data sets. For example, the

1 training data ( $f_i, t_i$ ) can be a hold out set (in which a fraction of the initial training  
2 set, such as 30%, is not used to train the SVM but is used to train the sigmoid) or  
3 can be generated using three-fold cross-validation (in which the initial training set  
4 is split into three parts, each of three SVMs is trained on permutations of two out  
5 of three parts, and the  $f_i$  are evaluated on the remaining third, and the union of all  
6 three sets  $f_i$  forming the training set of the sigmoid).

7 Additionally, an out-of-sample model is used to avoid "overfitting" the  
8 sigmoid. Out-of-sample data is modeled with the same empirical density as the  
9 sigmoid training data, but with a finite probability of opposite label. In other  
10 words, when a positive example is observed at a value  $f_i$ , rather than using  $t_i=1$ , it  
11 is assumed that there is a finite chance of opposite label at the same  $f_i$  in the out-  
12 of-sample data. Therefore, a value of  $t_i=1-\epsilon_+$  is used, for some  $\epsilon_+$ . Similarly, a  
13 negative example will use a target value of  $t_i=\epsilon_-$ .

14 Regardless of the manner in which the *posterior* probability  $P(C_i | X)$  for a  
15 0.5 second window is determined, the *posterior* probabilities for multiple windows  
16 are combined to determine the *posterior* probability for a segment. In one  
17 implementation, each segment is five seconds, so the *posterior* probabilities of ten  
18 adjacent windows are used to determine the *posterior* probability for each  
19 segment.

20 The *posterior* probabilities for the multiple windows can be combined in a  
21 variety of different manners. In one implementation, the *posterior* probability of  
22 the segment being an exciting segment, referred to as  $P(ES)$ , is determined by  
23 averaging the *posterior* probabilities of the windows in the segment:

$$24 \quad P(ES) = \frac{1}{M} \sum_{m=1}^M P(C_1 | X_m)$$

25



1 where  $C_l$  represents the excited speech class and  $M$  is the number of windows in  
2 the segment.

3 Which ten adjacent windows to use for a segment can be determined in a  
4 wide variety of different manners. In one implementation, if ten or more adjacent  
5 windows include speech, then those adjacent windows are combined into a single  
6 segment (e.g., which may be greater than ten windows, or, if too large, which may  
7 be pared down into multiple smaller ten-window segments). However, if there are  
8 fewer than ten adjacent windows, then additional windows are added (before  
9 and/or after the adjacent windows, between multiple groups of adjacent windows,  
10 etc.) to get the full ten windows, with the *posterior* probability for each of these  
11 additional windows being zero.

12 The probabilities  $P(ES)$  of these segments including excited speech 276 (as  
13 well as an indication of where these segments occur in the raw audio clip 250) are  
14 then made available to probabilistic fusion subsystem 264. Subsystem 264  
15 combines the probabilities 276 with information received from baseball hit  
16 detection subsystem 262, as discussed in more detail below.

17 Baseball hit detection subsystem 262 uses energy features 278 from feature  
18 extractor 252 to identify baseball hits within the audio data 250. In one  
19 implementation, the energy features 278 include the  $E_{23}$  and  $E_4$  features discussed  
20 above. Two additional features are also generated, which may be generated by  
21 feature extractor 252 or alternatively another component (not shown). These  
22 additional features are referred to as  $ER_{23}$  and  $ER_4$ , and are discussed in more  
23 detail below.

24 Hit detection is performed by subsystem 262 based on 25-frame groupings.  
25 A sliding selection of 25 consecutive 10ms frames of the audio data 250 is

analyzed, with the frame selection sliding frame-by-frame through the audio data 250. The features of the 25-frame groupings and a set of hit templates 280 are input to template matcher 282. Template matcher 282 compares the features of each 25-frame grouping to the hit templates 280, and based on this comparison determines a probability as to whether the particular 25-frame grouping contains a hit. An identification of the 25-frame groupings (e.g., the first frame in the grouping) and their corresponding probabilities are output by template matcher 282 as hit candidates 284.

Multiple-frame groupings are used to identify hits because the sound of a baseball hit is typically longer in duration than a single frame (which is, for example, only 10 ms). The baseball hit templates 280 are established to capture the shape of the energy curves (using the four energy features discussed above) over the time of the groupings (e.g., 25 10ms frames, or 0.25 seconds). Baseball hit templates 280 are designed so that the hit peak (the energy peak) is at the 8<sup>th</sup> frame of the 25-frame grouping. The additional features  $ER_{23}$  and  $ER_4$  are calculated by normalizing the  $E_{23}$  and  $E_4$  features based on the energy features in the 8<sup>th</sup> frame as follows:

$$ER_{23}(i) = \frac{E_{23}(i)}{E_{23}(8)}$$

$$ER_4(i) = \frac{E_4(i)}{E_4(8)}$$

where  $i$  ranges from 1 to 25,  $E_{23}(8)$  is the  $E_{23}$  energy in the 8<sup>th</sup> frame, and  $E_4(8)$  is the  $E_4$  energy in the 8<sup>th</sup> frame.

Fig. 6 illustrates exemplary baseball hit templates 280 that may be used in accordance with certain embodiments of the invention. The templates 280 in Fig. 6 illustrate the shape of the energy curves over time (25 frames) for each of the four features  $E_{23}$ ,  $E_4$ ,  $ER_{23}$ , and  $ER_4$ .

For each group of frames, template matcher 282 determines the probability that the group contains a baseball hit. This can be accomplished in multiple different manners, such as un-directional or directional template mapping. Initially, the four feature vectors for each of the 25 frames are concatenated, resulting in a 100-element vector. The templates 280 are similarly concatenated for each of the 25 frames, also resulting in a 100-element vector. The probability of a baseball hit in a grouping  $P(HT)$  can be calculated based on the Mahalanobis distance  $D$  between the concatenated feature vector and the concatenated template vector as follows:

$$D^2 = (\vec{X} - \vec{T})^T \Sigma^{-1} (\vec{X} - \vec{T})$$

where  $\vec{X}$  is the concatenated feature vector,  $\vec{T}$  is the concatenated template vector, and  $\Sigma$  is the covariance matrix of  $\vec{T}$ . Additionally,  $\Sigma$  is restricted to being a diagonal matrix, allowing the baseball hit probability  $P(HT)$  to be determined as follows:

$$P(HT) = \frac{\exp(-\frac{1}{2} D^2)}{C + \exp(-\frac{1}{2} D^2)}$$

1 where  $C$  is a constant that is data dependent (e.g.,  $\exp(-0.5D'^2)$ , where  $D'^2$  is the  
2 distance between the concatenated feature vector and a template for non-hit  
3 signals).

4 Alternatively, a directional template matching approach can be used, with  
5 the distance  $D$  being calculated as follows:

$$D^2 = (\vec{X} - \vec{T})^T I \times \Sigma^{-1} (\vec{X} - \vec{T})$$

6  
7  
8 where  $I$  is a diagonal indicator matrix. The indicator matrix  $I$  is adjusted to  
9 account for over-mismatches or under-mismatches (an over-mismatch is actually  
10 good). In one implementation, when the values of  $E_{23}$  for the 25-frame grouping  
11 are overmatching the templates (e.g., more than a certain number (such as one-  
12 half) of the data values in the 25-frame grouping are higher than the corresponding  
13 template values), then  $I = \text{diag}[1, \dots, 1, -1, 1, \dots, 1]$  where the  $-1$  is at location 8.  
14 However, when the values of  $E_{23}$  for the 25-frame grouping are under-matching  
15 the templates (e.g., less than a certain number (such as one-half) of the data values  
16 in the 25-frame grouping are less than the corresponding template values), then  $I =$   
17  $\text{diag}[-1, \dots, -1, -1, -1, \dots, -1]$  where the  $1$  is at location 8.

18 Although hit detection is described as being performed across all of the  
19 audio data 250, alternatively hit detection may be performed on only selected  
20 portions of the audio data 250. By way of example, hit detection may only be  
21 performed on the portions of audio data 250 that are excited speech segments (or  
22 speech windows) and for a period of time (e.g., five seconds) prior to those excited  
23 speech segments (or speech windows).  
24  
25

Probabilistic fusion generator 286 of subsystem 264 receives the excited speech segment probabilities  $P(ES)$  from excited speech classification subsystem 260 and the baseball hit probabilities  $P(HT)$  from baseball hit detection subsystem 262 and combines those probabilities to identify probabilities  $P(E)$  that segments of the audio data 250 are exciting. Probabilistic fusion generator 286 searches for hit frames within the 5-second interval of the excited speech segment. This combining is also referred to herein as "fusion".

Two different types of fusion can be used: weighted fusion and conditional fusion. Weighted fusion applies weights to each of the probabilities  $P(ES)$  and  $P(HT)$  adds the results to obtain the value  $P(E)$  as follows:

$$P(E) = W_{ES}P(ES) + W_{HT}P(HT)$$

where the weights  $W_{ES}$  and  $W_{HT}$  sum up to 1.0. In one implementation,  $W_{ES}$  is 0.83 and  $W_{HT}$  is 0.17, although other weights could alternatively be used.

Conditional fusion, on the other hand, accounts for the detected baseball hits adjusting the confidence level of the  $P(ES)$  estimation (e.g., that the excited speech probability is not high due to mislabeling a car horn as speech). The conditional fusion is calculated as follows:

$$\begin{aligned} P(E) &= P(CF)P(ES) \\ P(CF) &= P(CF | HT)P(HT) + P(CF | \overline{HT})P(\overline{HT}) \\ P(\overline{HT}) &= 1 - P(HT) \end{aligned}$$

where  $P(CF)$  is the probability of how much confidence there is in the  $P(ES)$  estimation, and  $P(\overline{HT})$  is the probability that there is no hit.  $P(CF|HT)$  represents the probability that we are confident that  $P(ES)$  is accurate given there is a



1 the pre-generated indications simply need to be accessed rather than determining,  
2 at the time of request, which segments to include in the summary.

3 Fig. 7 is a flowchart illustrating an exemplary process for rendering a  
4 program summary to a user in accordance with certain embodiments of the  
5 invention. The acts of Fig. 7 may be implemented in software, and may be carried  
6 out by a receiver 106 of Fig. 1 or alternatively a programming source of Fig. 1  
7 (e.g., Internet provider 120).

8 Initially, the user request for a summary is received along with parameters  
9 for the summary (act 300). The parameters of the summary identify what level of  
10 summary the user desires, and can vary by implementation. By way of example, a  
11 user may indicate as the summary parameters that he or she wants to be presented  
12 with any segments that have a probability of 0.75 or higher of being exciting  
13 segments. By way of another example, a user may indicate as the summary  
14 parameters that he or she wants to be presented with a 20-minute summary of the  
15 program.

16 The meta data corresponding to the program (the exciting segment  
17 probabilities  $P(E)$ ) is then accessed (act 302), and the appropriate exciting  
18 segments identified based on the summary parameters (act 304). Once the  
19 appropriate exciting segments are identified, they are rendered to the user (act  
20 306). The manner in which the appropriate exciting segments are identified can  
21 vary, in part based on the nature of the summary parameters. If the summary  
22 parameters indicate that all segments with a  $P(E)$  of 0.75 or higher should be  
23 presented, then all segments with a  $P(E)$  of 0.75 or greater are identified. If the  
24 summary parameters indicate that a 20-minute summary should be generated, then  
25 the appropriate segments are identified by determining (based on the  $P(E)$  of the

1 segments and the lengths of the segments) the segments having the highest  $P(E)$   
2 that have a combined length less than 20 minutes.

### 3 4 **Conclusion**

5 Although the description above uses language that is specific to structural  
6 features and/or methodological acts, it is to be understood that the invention  
7 defined in the appended claims is not limited to the specific features or acts  
8 described. Rather, the specific features and acts are disclosed as exemplary forms  
9 of implementing the invention.  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25



1 **CLAIMS**

2

3 1. One or more computer readable media having stored thereon a

4 plurality of instructions that, when executed by one or more processors, causes the

5 one or more processors to perform acts including:

6 receiving the audio portion of a sporting event;

7 classifying a set of segments of the audio portion as excited speech;

8 classifying a set of frame groupings as including baseball hits;

9 combining the set of segments and the set of frame groupings to identify

10 probabilities for each segment that the segment is an exciting segment; and

11 saving an indication of the set of segments and the corresponding

12 probabilities as meta data corresponding to the sporting event.

13

14 2. One or more computer readable media as recited in claim 1, wherein

15 each segment includes at least ten 0.5-second windows of the audio portion.

16

17 3. One or more computer readable media as recited in claim 1, wherein

18 each frame grouping is a grouping of 25 10-millisecond frames of the audio

19 portion.

20

21 4. One or more computer readable media as recited in claim 1, wherein

22 the classifying a set of segments as excited speech comprises:

23 extracting a first set of features from the audio portion;

24 identifying, based on the first set of features, a plurality of windows of the

25 audio portion that include speech;

1 extracting a second set of features from the audio portion; and  
2 identifying, based on the second set of features, which of the plurality of  
3 windows include excited speech.

4  
5 5. One or more computer readable media as recited in claim 4, wherein  
6 extracting the first set of features comprises, for each of a plurality of windows:

7 identifying an average waveform amplitude of the audio portion in a first  
8 frequency band of the window;

9 identifying an average waveform amplitude of the audio portion in a second  
10 frequency band of the window;

11 concatenating the identified average waveform amplitudes to generate an  
12 energy feature of the first set of features; and

13 determining, for an MFCC feature of the first set of features, the Mel-  
14 frequency Cepstral coefficient of the window.

15  
16 6. One or more computer readable media as recited in claim 5, wherein  
17 the identifying a plurality of windows of the audio portion that include speech  
18 comprises determining that a window includes speech if the energy feature  
19 exceeds a first threshold and the MFCC feature exceeds a second threshold.

20  
21 7. One or more computer readable media as recited in claim 4, wherein  
22 extracting the second set of features comprises, for each of a plurality of windows:

23 identifying, for each of a plurality of frames in the window, an average  
24 waveform amplitude of the audio portion in a first frequency band;

1 identifying, for each of the plurality of frames in the window, an average  
2 waveform amplitude of the audio portion in a second frequency band;  
3 concatenating the identified average waveform amplitudes to generate an  
4 energy feature;  
5 extracting, as a pitch feature of each frame, the pitch of each frame;  
6 identifying a plurality of statistics regarding each window based on the  
7 energy features and pitch features of the plurality of frames.

8  
9 8. One or more computer readable media as recited in claim 7, wherein  
10 the identifying a plurality of statistics further comprises:

11 identifying a maximum energy;  
12 identifying an average energy;  
13 identifying an energy dynamic range;  
14 identifying a maximum pitch;  
15 identifying an average pitch; and  
16 identifying a dynamic pitch range.

17  
18 9. One or more computer readable media as recited in claim 7, wherein  
19 the identifying which of the plurality of windows include excited speech  
20 comprises identifying a *posterior* probability that the window corresponds to an  
21 excited speech class and identifying a *posterior* probability that the window  
22 corresponds to a non-excited speech class, and classifying the window in the class  
23 with the higher *posterior* probability.  
24  
25

**10.** One or more computer readable media as recited in claim 4, further comprising instructions that cause the one or more processors to perform acts including outputting an excited speech probability for each of the segments that include excited speech, the excited speech probability for a segment indicating a likelihood that the segment includes excited speech.

11. One or more computer readable media as recited in claim 1, wherein the classifying a set of frame groupings as including baseball hits:

extracting, for each frame in a multiple-frame grouping, a set of features from the audio portion;

comparing the set of features from the multiple-frame groupings to a set of templates; and

identifying, for each of the multiple-frame groupings, a probability that the grouping includes a baseball hit based on how well the grouping matches the set of templates.

12. One or more computer readable media as recited in claim 11,  
wherein the extracting comprises, for each frame:

identifying an average waveform amplitude of the audio portion in a first frequency band of the frame;

identifying an average waveform amplitude of the audio portion in a second frequency band of the frame;

concatenating the identified average waveform amplitudes to generate a first energy feature of the set of features;



1           **16.**    A method comprising:  
2           receiving a program including both audio and video;  
3           receiving meta data corresponding to the program; and  
4           rendering, based on the meta data, portions of the program as a summary of  
5 the program.

6  
7           **17.**    A method as recited in claim 16, wherein the rendering comprises  
8 displaying the video of the portions and playing the audio of the portions.

9  
10          **18.**    A method as recited in claim 16, wherein the meta data comprises a  
11 probability indicator, for each of a plurality of portions of the program, that  
12 identifies a probability that the corresponding portion is an exciting portion of the  
13 program.

14  
15          **19.**    A method as recited in claim 18, wherein the rendering comprises  
16 selecting the plurality of portions that have probability indicators that exceed a  
17 threshold value, and rendering the selected portions as the summary.

18  
19          **20.**    A method as recited in claim 16, wherein the receiving a program  
20 and the receiving meta data comprise receiving both the program and the meta  
21 data from the same source.

22  
23          **21.**    A method as recited in claim 16, wherein the receiving meta data  
24 comprises receiving the meta data from a remote source via a network.  
25

1           **22.**    A method as recited in claim 21, wherein the network comprises the  
2 Internet.

3  
4           **23.**    A method as recited in claim 16, wherein the receiving meta data  
5 comprises receiving meta data generated manually.

6  
7           **24.**    A method as recited in claim 16, wherein the receiving meta data  
8 comprises receiving meta data generated automatically.

9  
10          **25.**    A method as recited in claim 16, wherein the meta data comprises a  
11 plurality of probabilities, each corresponding to a segment of the program, the  
12 probabilities representing a probabilistic combination of sports-specific events and  
13 sports-generic events identified in the program.

14  
15          **26.**    A method as recited in claim 25, wherein the sports-specific events  
16 comprise baseball hits, and wherein the sports-generic events comprise excited  
17 speech.

18  
19          **27.**    One or more computer readable media including a computer  
20 program that is executable by a processor to perform the method recited in claim  
21 16.

22          **28.**    A system comprising:  
23           a content provider to make programming content available to requesting  
24 clients;  
25





1           **33.**    A system as recited in claim 32, wherein the plurality of receivers  
2 are further to select, from the plurality of portions, portions that have probability  
3 indicators that exceed a threshold value, and render the selected portions as the  
4 summary.

5  
6           **34.**    A system as recited in claim 28, wherein the meta data is generated  
7 manually.

8  
9           **35.**    A system as recited in claim 28, wherein the meta data is generated  
10 automatically.

11  
12           **36.**    A system as recited in claim 28, wherein the meta data comprises a  
13 plurality of probabilities, each corresponding to a segment of the programming  
14 content, the probabilities representing a probabilistic combination of sports-  
15 specific events and sports-generic events identified in the programming content.

16  
17           **37.**    A system as recited in claim 36, wherein the sports-specific events  
18 comprise baseball hits, and wherein the sports-generic events comprise excited  
19 speech.

20  
21           **38.**    A method of automatically summarizing a program, the method  
22 comprising:

23           identifying a plurality of sports-generic events from the audio of the  
24 program;  
25

1 identifying a plurality of sports-specific events from the audio of the  
2 program; and

3 identifying, by combining the sports-generic events and the sports-specific  
4 events, a set of portions of the program as a summary of the program.

5  
6 **39.** A method as recited in claim 38, further comprising transmitting the  
7 set of portions to a client computer as the summary of the program.

8  
9 **40.** A method as recited in claim 39, wherein the transmitting comprises  
10 transmitting the set of portions via the Internet.

11  
12 **41.** A method as recited in claim 38, wherein the sports-specific events  
13 comprise baseball hits, and wherein the sports-generic events comprise excited  
14 speech.

15  
16 **42.** A method as recited in claim 38, wherein the program includes both  
17 an audio portion and a video portion.

18  
19 **43.** One or more computer readable media including a computer  
20 program that is executable by a processor to perform the method recited in claim  
21 38.

22  
23 **44.** A method comprising:  
24 analyzing audio data of a program to identify a first plurality of portions of  
25 the program each including excited speech;

1 analyzing the audio data to identify a second plurality of portions of the  
2 program each including a potential baseball hit; and

3 combining the first plurality of portions and the second plurality of portions  
4 to identify a set of segments of the program and a likelihood, for each of the  
5 segments in the set, that the segment is an exciting part of the program.

6  
7 **45.** A method as recited in claim 44, wherein the analyzing audio data to  
8 identify the first plurality of portions comprises:

9 extracting a first set of features from the audio data;

10 identifying, based on the first set of features, a plurality of windows of the  
11 audio data that include speech;

12 extracting a second set of features from the audio data; and

13 identifying, based on the second set of features, which of the plurality of  
14 windows include excited speech.

15  
16 **46.** A method as recited in claim 45, wherein extracting the first set of  
17 features comprises, for each of a plurality of windows:

18 identifying an average waveform amplitude of the audio data in a first  
19 frequency band of the window;

20 identifying an average waveform amplitude of the audio data in a second  
21 frequency band of the window;

22 concatenating the identified average waveform amplitudes to generate an  
23 energy feature of the first set of features; and

24 determining, for an MFCC feature of the first set of features, the Mel-  
25 frequency Cepstral coefficient of the window.

1  
2       **47.**    A method as recited in claim 46, wherein each window comprises  
3 0.5 seconds.

4  
5       **48.**    A method as recited in claim 46, wherein the identifying a plurality  
6 of windows of the audio data that include speech comprises determining that a  
7 window includes speech if the energy feature exceeds a first threshold and the  
8 MFCC feature exceeds a second threshold.

9  
10       **49.**    A method as recited in claim 45, wherein extracting the second set  
11 of features comprises, for each of a plurality of windows:

12       identifying, for each of a plurality of frames in the window, an average  
13 waveform amplitude of the audio data in a first frequency band;

14       identifying, for each of the plurality of frames in the window, an average  
15 waveform amplitude of the audio data in a second frequency band;

16       concatenating the identified average waveform amplitudes to generate an  
17 energy feature;

18       extracting, as a pitch feature of each frame, the pitch of each frame; and

19       identifying a plurality of statistics regarding each window based on the  
20 energy features and pitch features of the plurality of frames.

21  
22       **50.**    A method as recited in claim 49, wherein the identifying a plurality  
23 of statistics further comprises:

24       identifying a maximum energy;

25       identifying an average energy;





1 identifying an average waveform amplitude of the audio data in a second  
2 frequency band of the frame;

3 concatenating the identified average waveform amplitudes to generate a  
4 first energy feature of the set of features;

5 identifying an average waveform amplitude of the audio data in a third  
6 frequency band of the frame; and

7 using, as a second energy feature of the set of features, the average  
8 waveform amplitude of the third frequency band.

9  
10 **59.** A method as recited in claim 58, further comprising:

11 generating a third energy feature for each frame by normalizing the first  
12 energy feature based on the first energy feature of the eighth frame of the multiple-  
13 frame grouping; and

14 generating a fourth energy feature for each frame by normalizing the  
15 second energy feature based on the second energy feature of the eighth frame of  
16 the multiple-frame grouping.

17  
18 **60.** A method as recited in claim 44, wherein the combining comprises  
19 generating, for each of the first plurality of portions, a weighted sum of a  
20 probability that the portion includes excited speech and a probability that the  
21 portion includes a baseball hit.

1           **61.**    A method as recited in claim 44, wherein the combining comprises  
2 adjusting the probability that a portion includes excited speech based on the  
3 probability that the portion includes a baseball hit.

4  
5           **62.**    A system comprising:  
6           a feature extractor to extract a plurality of audio features from programming  
7 content;  
8           an excited speech classification subsystem to identify, based on a sub-set of  
9 the audio features, a set of segments of the programming content and  
10 corresponding probabilities that the segments include excited speech;  
11           a baseball hit detection subsystem to identify, based on another sub-set of  
12 the audio features, a set of frame groupings of the programming content and  
13 corresponding probabilities that the frame groupings include baseball hits; and  
14           a probabilistic fusion subsystem to combine the probabilities that the  
15 segments include excited speech and the probabilities that the frame groupings  
16 include baseball hits, and to generate a probability that portions of the  
17 programming content are exciting based on the combination.

18  
19           **63.**    A system as recited in claim 62, wherein the excited speech  
20 classification subsystem is to identify the set of segments by:

21           identifying, based on a first set of the plurality of audio features, a plurality  
22 of windows of the programming content that include speech;

23           identifying, based on a second set of the plurality of audio features, which  
24 of the plurality of windows include excited speech.



1           **64.**    A system as recited in claim 62, wherein the baseball hit detection  
2 subsystem is to identify the set of frame groupings by:

3           combining, for each of a plurality of multiple-frame groupings, a set of  
4 features from the programming content;

5           comparing the sets of features from the multiple-frame groupings to a set of  
6 templates; and

7           identifying, for each of the multiple-frame groupings, a probability that the  
8 grouping includes a baseball hit based on how well the grouping matches the set of  
9 templates.

10  
11           **65.**    A system comprising:

12           a receiving device to receive a sporting event;

13           a user interface to receive, from a user, an indication of a desired summary  
14 level for the sporting event; and

15           a processing subsystem to identify which portions of the sporting event to  
16 present to the user based at least in part on both the desired summary level and  
17 meta data corresponding to the sporting event, the meta data identifying a  
18 likelihood of each of a plurality of portions of the sporting event being exciting  
19 based at least in part on the presence of both excited speech and ball hits within  
20 the sporting event.

1 **ABSTRACT**

2 Audio/video programming content is made available to a receiver from a  
3 content provider, and meta data is made available to the receiver from a meta data  
4 provider. The meta data corresponds to the programming content, and identifies,  
5 for each of multiple portions of the programming content, an indicator of a  
6 likelihood that the portion is an exciting portion of the content. In one  
7 implementation, the meta data includes probabilities that segments of a baseball  
8 program are exciting, and is generated by analyzing the audio data of the baseball  
9 program for both excited speech and baseball hits. The meta data can then be used  
10 to generate a summary for the baseball program.

11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

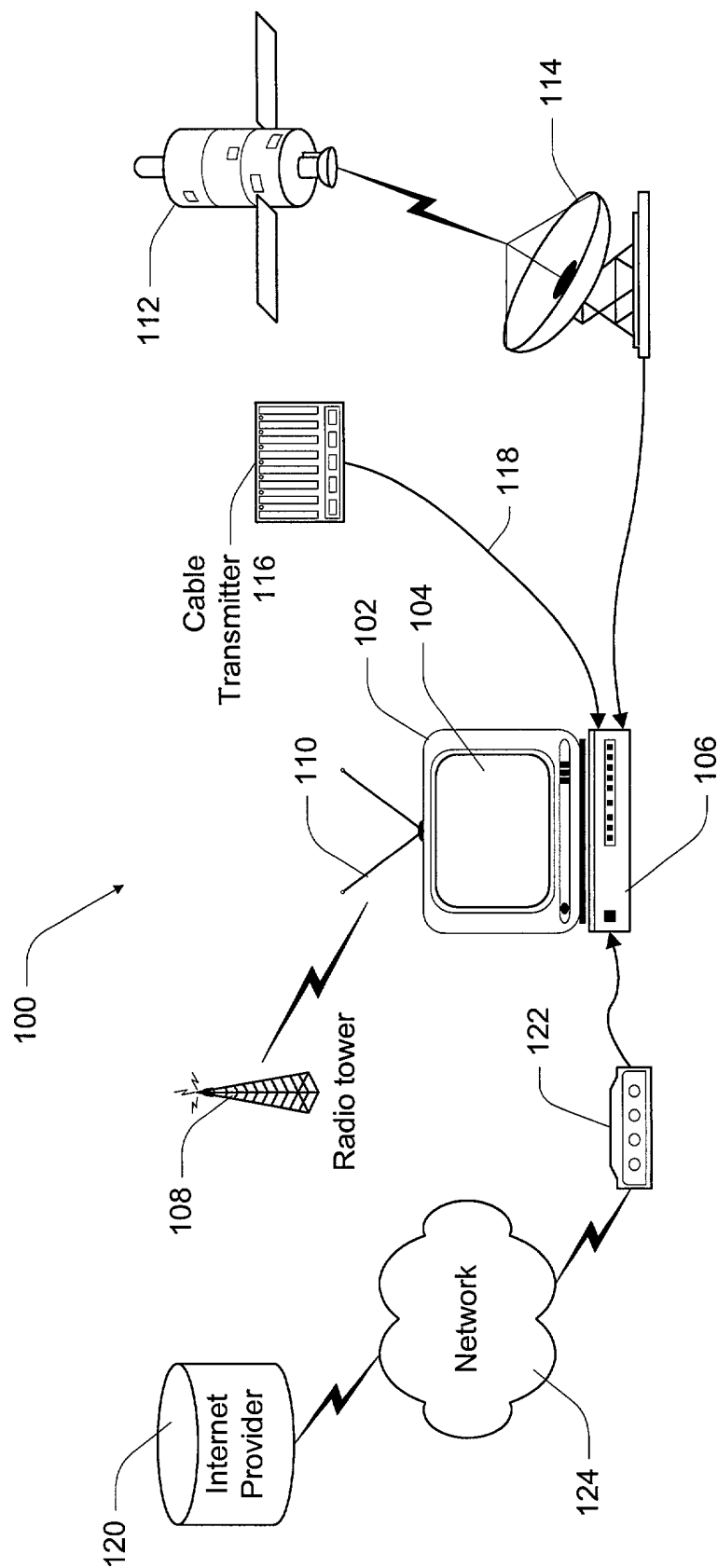
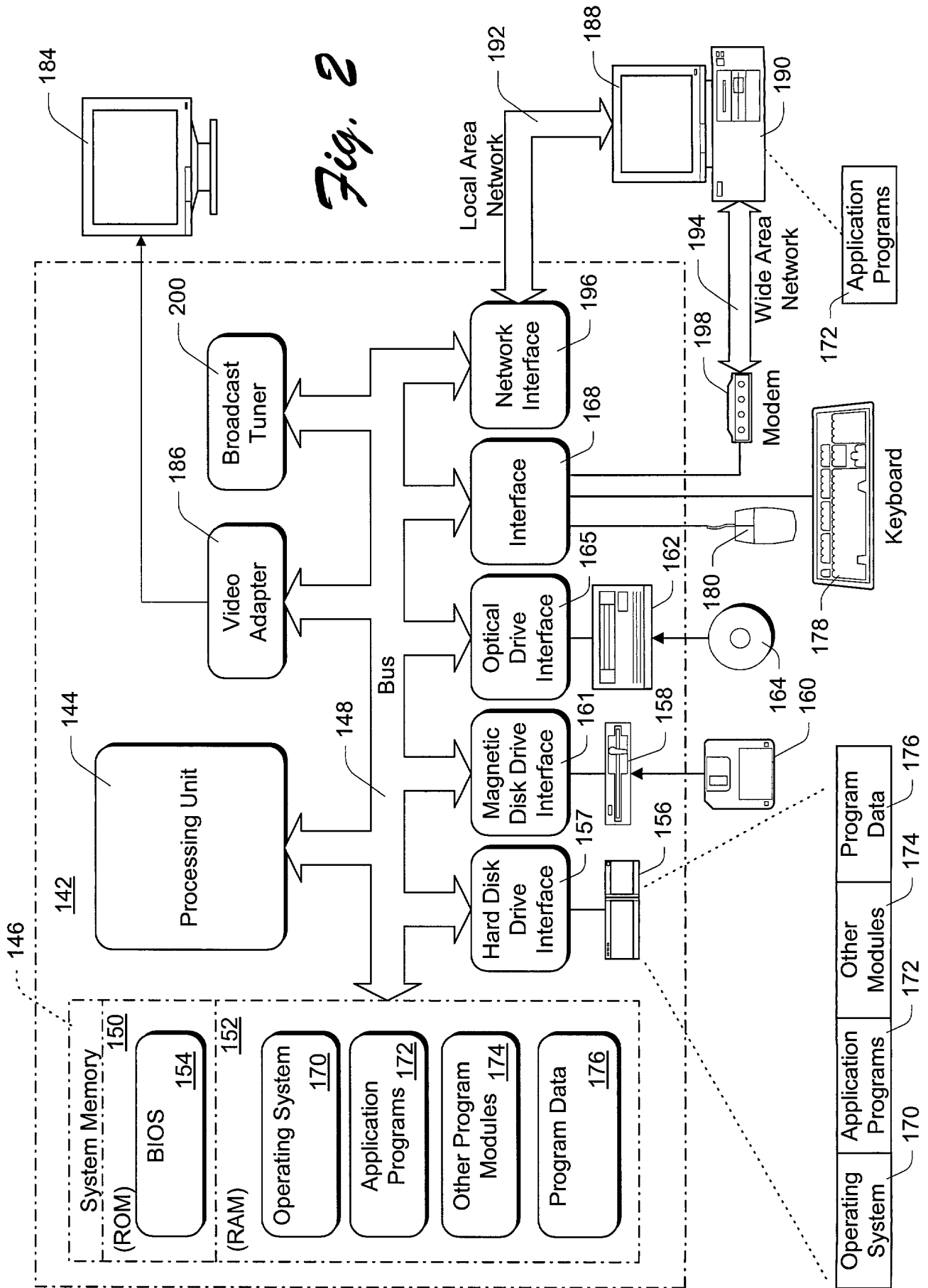


Fig. 1

Fig. 2



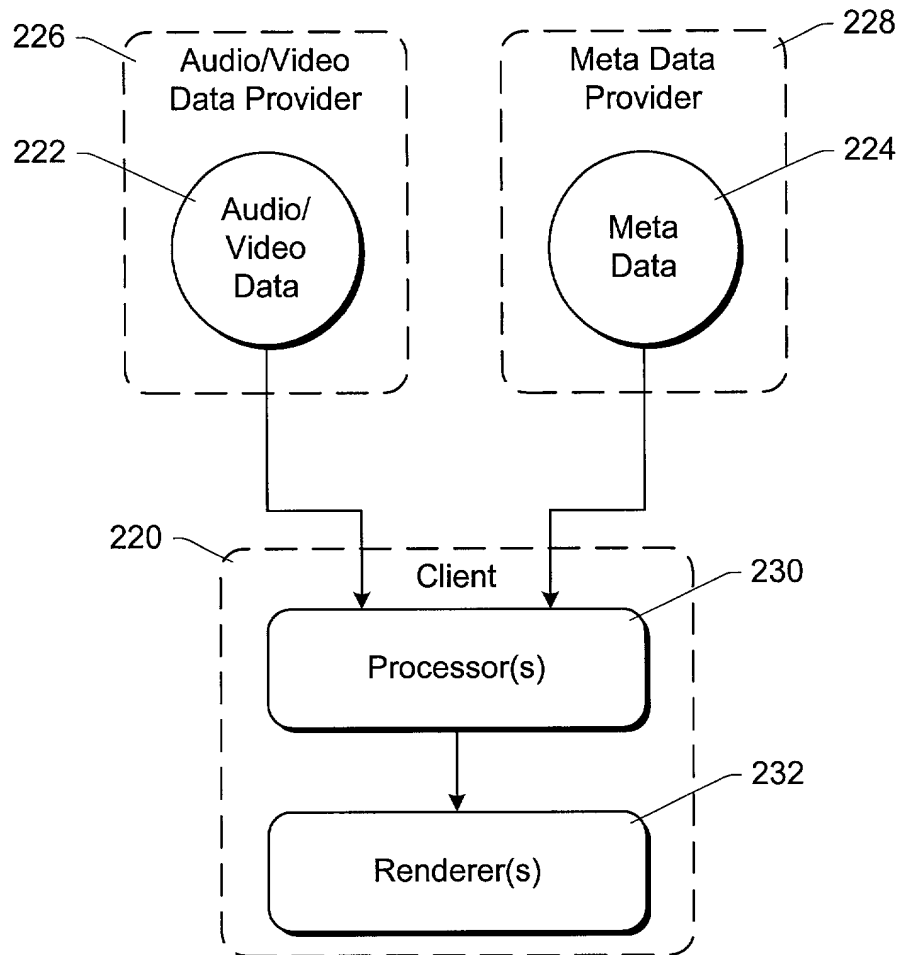
*Fig. 3*



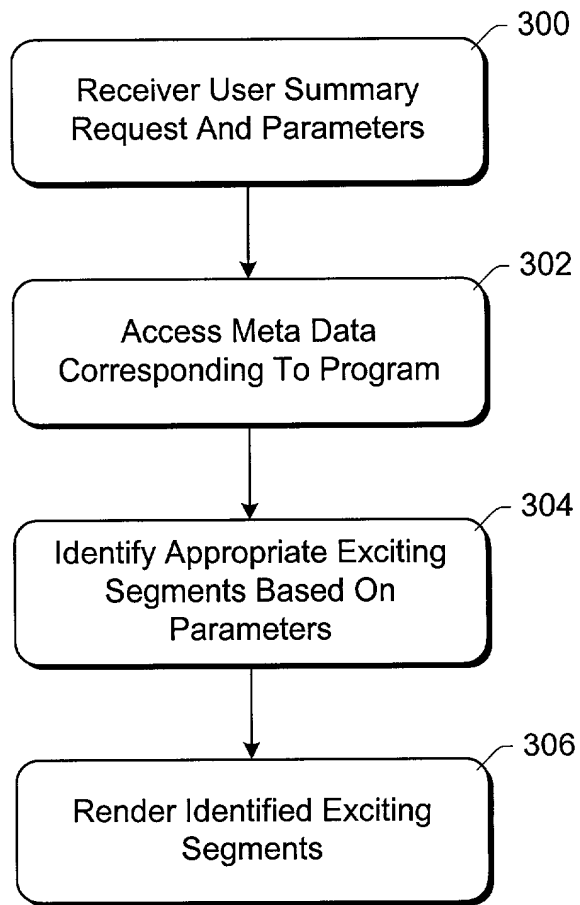


Fig. 5



*Fig. 6*





*Fig. 7*